**IMG 2.4.1** released on **September 1$^{st}$, 2008** is an updated version of IMG/M 2.4 which was released in February 2008.

# IMG/M 2.4.1 Content

**IMG/M 2.4** (released in Feb 2008) contains **reference genomes** from **IMG 2.4** (released in Dec 2007) integrated with **metagenome** datasets generated from samples for the following metagenome projects:

- A methane-oxidizing archaeal community sample
- two biological phosphorus removing (EBPR) sludge samples;
- three isolated deep sea "whale fall" carcasses;
- an agricultural soil sample;
- an acid mine drainage (AMD) biofilm sample;
- two human distal gut samples;
- four gutless marine worm sample;
- five obese and lean mouse gut samples;
- three single cell TM7 samples;
- two termite hindgut samples;
- a uranium contaminated groundwater samples.

**IMG/M 2.4.1,** released in Sep 2008, has been updated with datasets generated from the following microbial community samples:

- ten hypersaline microbial mat samples[1],
- five samples from Lake Washington in Seattle[2], and
- two airborne samples from in an indoor urban environment[3].

---

[1] http://www.nature.com/msb/journal/v4/n1/pdf/msb200835.pdf
[2] http://www.nature.com/nbt/journal/vaop/ncurrent/abs/nbt.1488.html
[3] http://www.plosone.org/article/info:doi/10.1371/journal.pone.0001862

# IMG/M Statistics

Various statistics are provided via the **IMG/M Statistics** link on the home page of IMG/M. The list of microbiome samples grouped by study or project is provided via the **Microbiome Samples** link, with a link to the associated publication also provided for each project. A **Map** link on the home page provides a Google Map showing the location of these samples, as shown below. For each sample, a link to its **Microbiome Details** page is provided.

# IMG/M User Interface

The User Interface has been extended in order to improve its overall usability.  The main extensions include **graphical viewers,** such as for **protein family coverage** of single or multiple genomes and metagenomes, and an **Abundance Profile Overview** tool that extends the Abundance Profile Viewer.

# Graphical Viewers

## Microbiome Details – *Microbiome Information, Metagenome Statistics*

The **Microbiome Information** part of **Microbiome Details** contains additional metadata collected from scientists as well as a map display of the sample collection location using Google Map.

The **Metagenome Statistics** part of **Microbiome Details** has been reorganized in order to improve clarity. Graphical viewers have been added for displaying the distribution of genes associated with COG, Pfam, TIGRfam, and KEGG, as illustrated below.
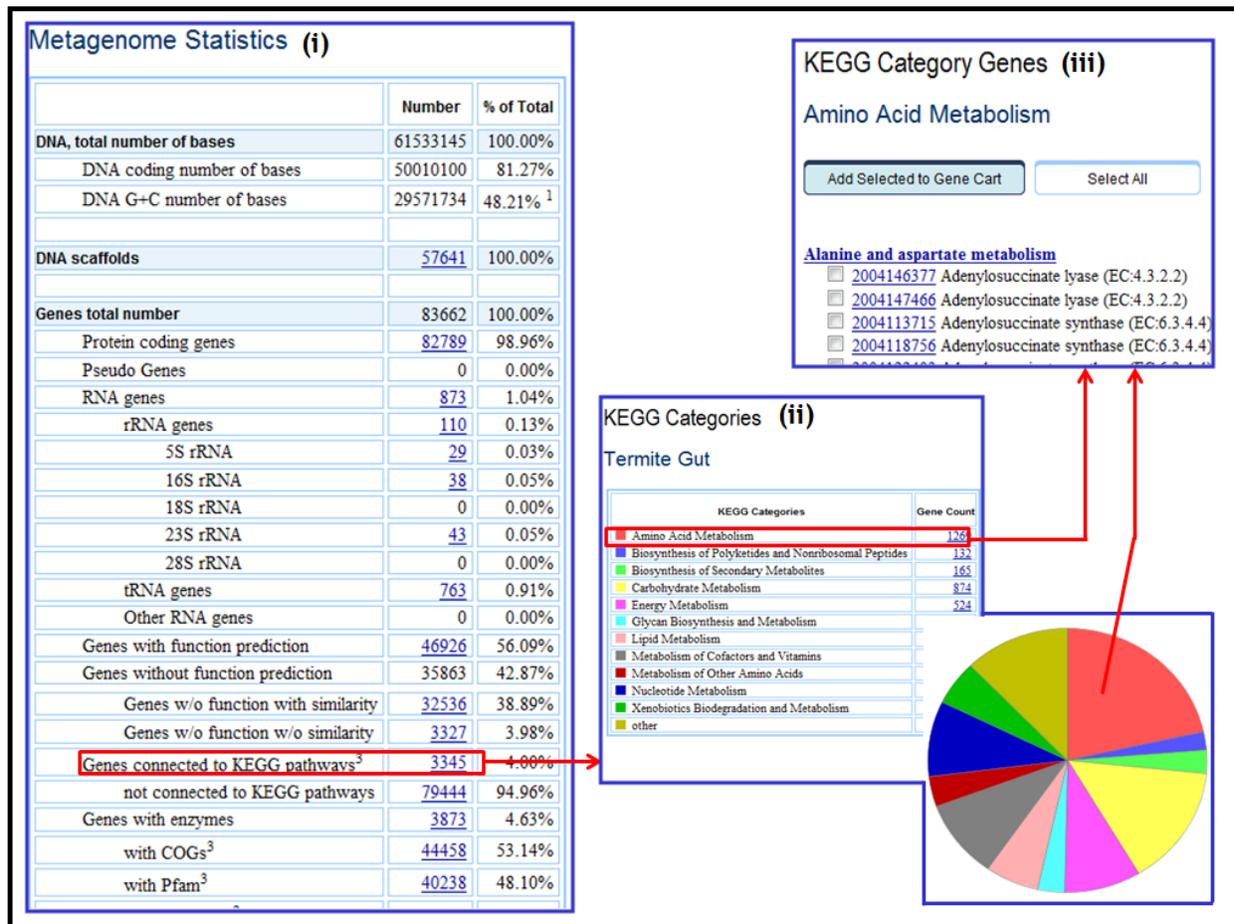


**FIGURE 1. Organism Details - Genome Statistics.**

Gene counts in some categories in **Metagenome Statistics**, such as "Genes with  COGs", "Genes with Pfam", "Genes with TIGRfam", and "Genes connected to KEGG pathways" are

linked to tables that show these genes classified according to the corresponding functional hierarchies (e.g., COG Functional Categories, KEGG Categories, etc.), as illustrated in Figure 1(ii), displayed both in tabular and graphical (pie chart) format.

The gene counts in the table and on the pie chart representing (e.g., COG, KEGG) functional categories are linked to the table that contains groupings of genes according to individual functional groups or metabolic pathways, as illustrated in Figure 1(iii). Genes can be then selected and saved in the **Gene Cart** for further analysis.

## Microbiome  Details – *Genome Viewers*

The **Genome Viewers** part of **Microbiome Details** includes a new display for scaffolds and an additional version of the **Chromosome Viewer**, as illustrated in Figure 2. Scaffolds are displayed in two lists sorted by gene count and sequence length, respectively. Each list is also displayed graphically using bar charts, as illustrated in Figure 2(ii).

After a scaffold and coordinate range is selected, as illustrated in Figure 2(iii), the **Chromosome Viewer** displays genes colored by their associated COG as illustrated in Figure 2(iv). Gene coloring can be then switched to reflect deviation of characteristic GC percentage for that genome, as illustrated in Figure 2(v).



**FIGURE 2. Microbiome Details - Genome Viewers.**

# Microbiome Details – *Phylogenetic Distribution of Genes*

The **Phylogenetic Distribuiton of Genes** part of **Microbiome Details** has been extended with graphical displays for viewing **COG Functional Category Statistics** and **COG Pathway Statistics,** as illustrated in Figure 3, in addition to the existing tabular displays.

Selecting the "Chart" option for viewing one of the **COG Functional Category Statistics** or **COG Pathway Statistics** (see Figure 3(ii)), will display in a tabular and *pie chart* format the count of genes across all phyla/classes associated with each COG category or pathway, as illustrated in Figure 3(iii). Clicking on a COG category or pathway on the pie chart or on the colored coded square for a COG category or pathway in the table, will display a *bar chart* with the number of genes for each phylum/class associated with that COG category or pathway, as illustrated on the lower side pane of Figure 3(iii).
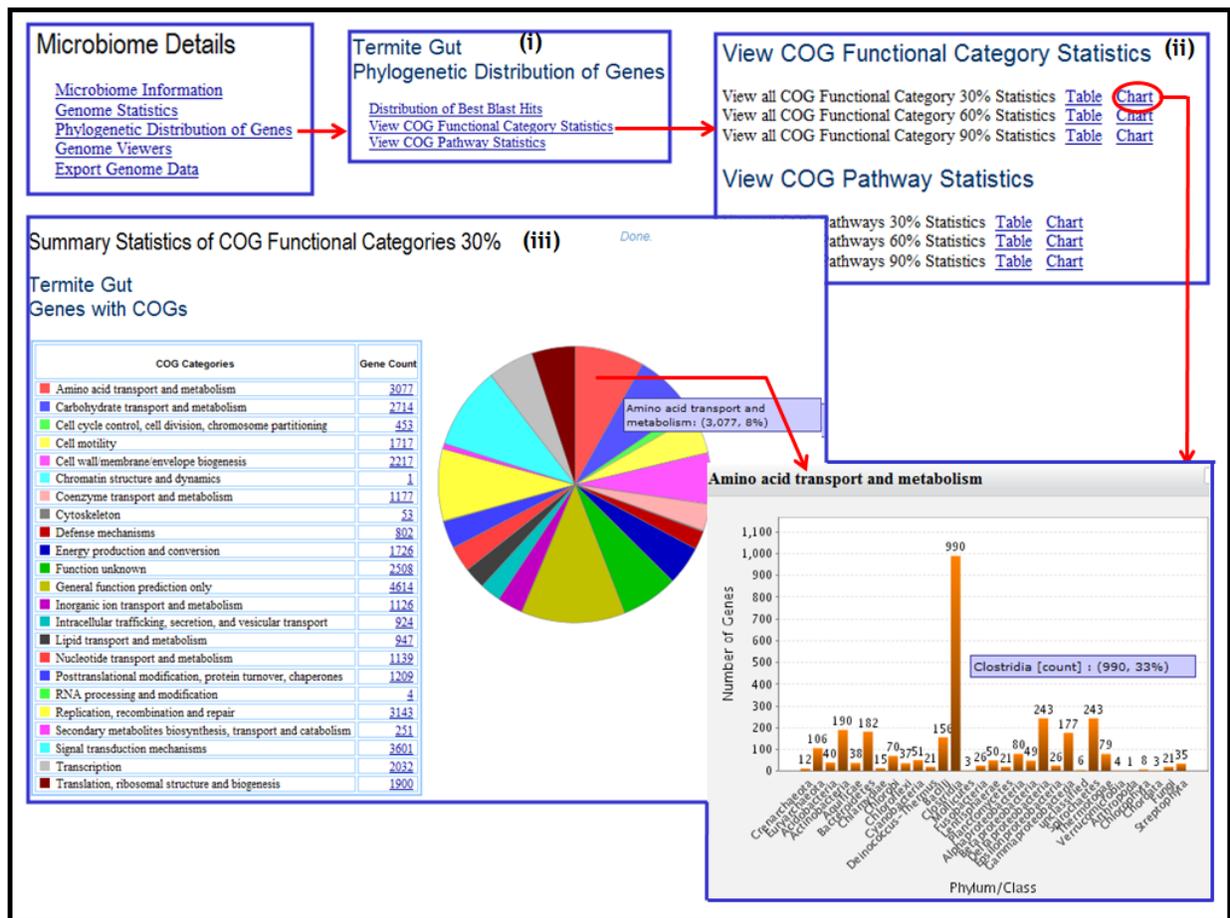


**FIGURE 3. Microbiome Details – Phylogenetic Distribution of Genes.**

# Gene Details – *Evidence for Function Prediction*

The **Evidence for Function Prediction** part of **Gene Details** includes a new **Sequence Viewer** and an additional version of the **Chromosome Viewer**, as illustrated in Figure 4(i).

For a specific gene, the **Sequence Viewer** displays the six frame translation with putative ORF's, potential start codons, and potential Shine-Delgano regions. The gene neighborhood, minimum ORF size and type of display (graphic or text) can be selected, as illustrated in Figure 4(ii).The text display provides the protein sequences for the ORFs while the graphical display, illustrated in Figure 4(iii), includes a GC plot. The additional **Chromosome Viewer** displays the neighborhood of the target gene with genes colored to reflect deviation of characteristic GC percentage for that genome, as illustrated in Figure 4(iv).



**FIGURE 4. Gene Details – Evidence for Function Prediction Viewers.**

## Compare Genomes – *Genome Statistics*

**Graphical viewers** have been added for displaying the results of **COG and KEGG Category Statistics**, as illustrated below for COG Category Statistics.

Selecting **Statistics for Genomes by specific COG Category** (see Figure 5(i)) will display in a tabular and pie chart format the count of genes associated with each COG category across all selected genomes, or metagenomes as illustrated in Figure 5(ii). Clicking on a COG category on the *pie chart* or on the colored coded square for a COG category in the table will display a *bar chart* with the percent of genes for each genome or metagenome associated with that COG category, as illustrated on the lower side pane of Figure 5(ii).
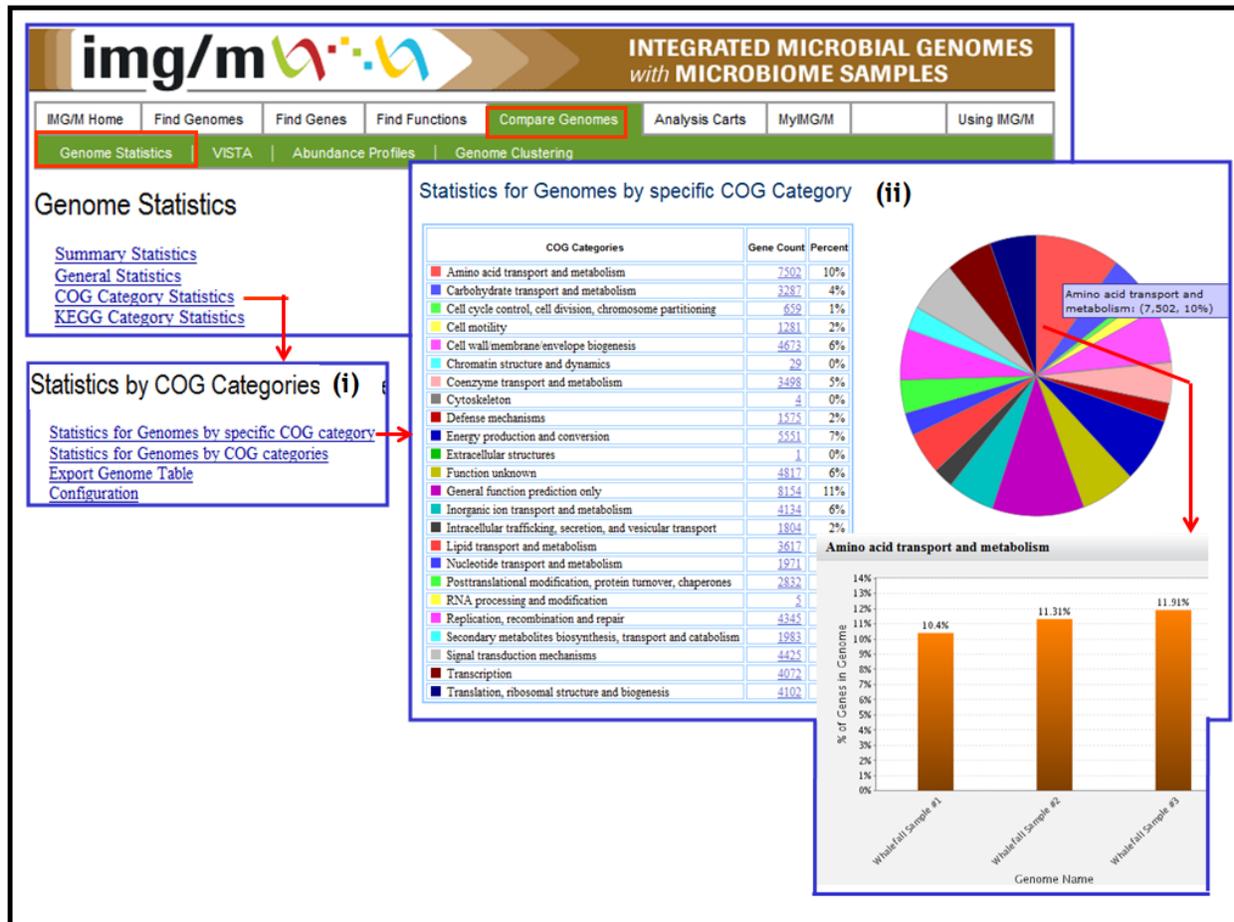


**FIGURE 5. Genome Statistics: COG Category statistics.**

# Abundance Profile Overview

The **Abundance Profile Overview** illustrated in Figure 6 extends the **Abundance Profile Viewer**. First, select the type of format for displaying the results ("Heat Map" or "Matrix"), protein/functional families (COG, Pfam, TIGRfam, Enzyme), normalization method, and a set of metagenomes or genomes in the **Abundance Profile Viewer** page, as illustrated in Figure 6(i) For "Heat Map" output, the abundance of protein/functional families is displayed as a heat map with red corresponding to the most abundant families, as illustrated in Figure 6(ii). Each column on the map corresponds to a genome or metagenome, and each row corresponds to a family. Click on the cell in order to retrieve the list of genes assigned to this particular family in a genome or metagenome. Click on the identifier of the family displayed on the right of the column in order to include the corresponding family into the **Function Cart**.



**FIGURE 6. Abundance Profile Overview: Results with Heat Map and Matrix Output Format.**

If the "Matrix" output is selected, the abundance of protein/functional families is displayed in a tabular format, with each row corresponding to a family and each cell containing the number of genes associated with a family for a specific genome or metagenome, as illustrated in Figure 6(iii). Click on the cell in order to retrieve the list of genes assigned to this particular family in a genome or metagenome. Families of interest can be selected for inclusion into the **Function Cart**. The results in "Matrix" format can be exported to a tab-delimited Excel file.

# Function Comparison

The **Function Comparison** analysis tool has a new organization of the function comparison results and an updated computation of the statistical significance.

The function comparison result lists for each function, **F**: (the number of genes or estimated gene copies in the query genome/metagenome, **Q**, associated with **F**, and for each reference genome/metagenome, **R**, the number of genes or estimated gene copies in **R** associated with **F**, together with the D-score and p-value associated with the comparison of **Q** to **R**, as illustrated in Figure 7.

In the computation of D-scores, a False Discovery Rate (FDR) correction[4] is used prior to the display of significant hypotheses.  FDR correction is meant to control the expected number of false predictions in a multiple-testing scenario. For more details, see section 3.3 of UsingIMG/M (http://img.jgi.doe.gov/m/doc/userGuide_m.pdf).



**FIGURE 7. Function Comparison.**

---

[4] Benjamini Y and Hochberg Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society* **57** (1), 289–300.